

ON SELECTING A SUBSET CONTAINING THE POPULATION

WITH THE SMALLEST VARIANCE ^{1/}

Shanti S. Gupta
Bell Telephone Laboratories

Milton Sobel
University of Minnesota

Special Report No. 4

Department of Statistics
University of Minnesota
Minneapolis, Minnesota

^{1/} This work was started and the tables were constructed at the Bell Telephone Laboratories.

ON SELECTING A SUBSET CONTAINING THE
POPULATION WITH THE SMALLEST VARIANCE

by

Shanti S. Gupta and Milton Sobel*
Bell Telephone Laboratories

1. Summary

A multiple decision approach is taken to the problem of selecting a subset from k given normal populations which includes the "best" population, i.e., the one with the smallest population variance. The population variances of the k normal populations are unknown and the population means may be known or unknown. Based on a common number of observations from each population, a procedure R is defined which selects a subset which is never empty, small in size and yet large enough to guarantee with preassigned probability that it includes the best population, regardless of what are the true unknown population variances. Expressions for the probability of a correct selection using R are derived and it is shown that, for the case in which the k sample variances have a common number, ν , of degrees of freedom, the infimum of this probability is identical with the probability integral of the ratio of the minimum of $k-1$ independent chi-squares to another independent chi-square, all with ν degrees of freedom. The associated distribution theory for this statistic and the tables needed to carry out the procedure R are given in a companion paper [6]. Formulas are obtained for the expected number of populations retained in the selected subset and it is shown that this function attains its maximum when the population variances are all equal. Two generalizations are considered; one deals with the case of unequal degrees of freedom and the other is concerned with a procedure for the selection of the t ($1 \leq t < k$) best populations.

*Now at the University of Minnesota

2. Introduction

As alternatives to the classical tests of homogeneity in the Analysis of Variance, new techniques have been developed in recent years which try to incorporate in the original statistical formulation of the problem the plans of the experimenter for further analysis after the hypothesis of homogeneity is tested. For example, if the experimenter tests for homogeneity and regardless of the outcome, ranks his populations on the basis of further analysis of the same data, it would be more realistic to assume at the outset that the parameters are unequal and formulate the main problem as a ranking problem.

The formulation considered in this and earlier papers [3], [4], [5] is that of selecting a subset of k populations which contains the "best" population where the best population is usually defined as the one with the largest (or smallest) parameter value; some further remarks on the motivation of this formulation are contained in these papers.

In the present paper the parameters of interest are the unknown variances of the k normal populations with all means known or all means unknown. The object is to select a subset which includes the population with the smallest variance with a preassigned probability P^* , regardless of the true values of the k variances. The procedure R (defined in Section 3) depends only on the sample variances each of which (properly normalized) has a chi-square distribution and hence, under the same formulation, the results of this paper can be applied to any set of k chi-square statistics with a common number of degrees of freedom.

It should be pointed out that this "selecting a subset" formulation is different from and in a certain respect related to the "indifference zone" formulation for the problem of ranking variances treated in [2]. In the latter formulation, an indifference zone in the parameter space is preassigned, the common number of observations needed is tabulated, and the final decision is the selection of a single population which is asserted to be the best population. In the formulation of this paper the number of observations is given, constants needed for carrying out the procedure are tabulated, and the final decision is the selection of a subset of populations which is asserted to contain the best population.

There are several aspects of this formulation which have already been described in the paper dealing with k binomial populations [5] which also apply to this problem. One is that R can be regarded as an elimination or screening procedure. Another is that a confidence statement can be made after experimentation. A third is that the expected size of the selected subset can be regarded as a measure of the efficiency of the procedure.

The main problem is formally described in Section 3 and the procedure R is defined in Section 4. In Section 5 we derive exact and asymptotic expressions for the probability of a correct selection using procedure R and the infimum of this probability over all points in the parameter space. Section 6 deals with the expected number of populations retained in the selected subset. Two generalizations are considered in Section 7, one dealing with the case of unequal degrees of freedom and the other describing a procedure for the problem of selecting a subset containing the t best populations.

3. Formal Statement of the Problem

Let $\Pi_1, \Pi_2, \dots, \Pi_k$ denote k given normal populations with unknown variances $\sigma_1^2, \sigma_2^2, \dots, \sigma_k^2$, respectively, (where each $\sigma_i^2 > 0$) and with all means known or all means unknown. The ordered variances are denoted by

$$(3.1) \quad \sigma_{[1]}^2 \leq \sigma_{[2]}^2 \leq \dots \leq \sigma_{[k]}^2,$$

(equalities being allowed for mathematical convenience only). It is assumed that there is no a priori information available about the correct pairing of the k given populations and the ordered scale parameters $\sigma_{[i]}^2$.

The population with variance equal to $\sigma_{[1]}^2$ is called the best population. The goal is to select a subset of the k populations containing the best population. Any such selection will be called a correct selection (CS). Then the problem is to find a rule R such that for a preassigned probability P^*

$$(3.2) \quad P \left\{ \text{CS} \mid R \right\} \geq P^*,$$

regardless of the true unknown values of the populations variances. It is assumed that the same number n of observations will be taken from each population; the case of unequal numbers of observations is briefly considered in Section 7.

From each population Π_i ($i = 1, 2, \dots, k$) we take n observations and, if the mean μ_i is known, we compute the statistic

$$(3.2) \quad s_i^2 = \frac{1}{n} \sum_{j=1}^n (x_{ij} - \mu_i)^2 \quad (i = 1, 2, \dots, k)$$

which, when multiplied by v/σ^2 (here $v = n$) has the χ_v^2 distribution with v degrees of freedom. [If the population means are unknown then we use

$$(3.3) \quad s_i^2 = \frac{1}{n-1} \sum_{j=1}^n (x_{ij} - \bar{x})^2$$

and the value of v is then $n-1$]. These statistics s_i^2 ($i=1, 2, \dots, k$) form a set of sufficient statistics for the problem and the rule R (defined in Section 4) depends only on these statistics.

It is clear that we prefer rules which make the size S of the selected subset never empty and as small as possible, subject to satisfying (3.2). [One can always attain any specified P^* , even unity, by putting all the populations in the selected subset.]

The rule R is defined in Section 4 for each positive integer $n \geq 2$; let it be denoted here by $R(n)$. Let $\vec{\sigma}^2$ denote the vector of true ordered variances which can vary in a space Ω , in which all vectors have only positive ordered components. Let $\Omega(\delta)$ be that part of Ω in which $\sigma_{[1]}^2 \leq \delta \sigma_{[2]}^2$ where $\delta > 0$ is preassigned. Finally, let $\epsilon > 0$ be preassigned. Two secondary problems can now be stated. One is to find the smallest common sample size n such that for some particular point $\vec{\sigma}_0^2$ in Ω

$$(3.4) \quad E \left\{ S; k, \vec{\sigma}_0^2, P^*, R(n) \right\} \leq 1 + \epsilon.$$

Another problem is to find the smallest common sample size n such that

$$(3.5) \quad \sup_{\vec{\sigma}^2 \text{ in } \Omega(\delta)} E \left\{ S; k, \vec{\sigma}^2, P^*, R(n) \right\} \leq 1 + \epsilon.$$

The expected size S of the retained subset is regarded as analogous with the complement of the "power" of the test of a hypothesis and both (3.4) and (3.5) are conditions which insure good "power". It is assumed here that both ϵ and δ (or both ϵ and σ_0^2) can be specified by the experimenter.

4. Procedure R

Let the ordered values of the k observed sample variances s_i^2 ($i = 1, 2, \dots, k$) all based on a common number ν of degrees of freedom be noted by

$$(4.1) \quad s_{[1]}^2 \leq s_{[2]}^2 \leq \dots \leq s_{[k]}^2.$$

The procedure R is then defined as follows. Procedure R: "Retain Π_1 in the selected subset if and only if

$$(4.2) \quad s_1^2 \leq s_{[1]}^2/c$$

where $c = c(\nu, k, P^*)$ is a constant with $0 < c \leq 1$ which is determined in advance of experimentation."

The constant c is chosen to be the largest value which satisfies the basic probability requirement (3.2) for all true configurations $\vec{\sigma}^2 = (\sigma_1^2, \sigma_2^2, \dots, \sigma_k^2)$. Tables of c -values for $\nu = 2(2)50$, $k = 2(1)11$ and $P^* = .75, .90, .95$, and $.99$ are given in a companion paper [6]; these c -values can also be regarded as percentage points of a smallest Studentized χ^2 -statistic as explained in Section 5.

Illustration

From each of $k=5$ populations Π_1 , a statistic T_1 is computed such that T_1/θ_1 has a χ_ν^2 distribution where the common value of ν is 14. It is desired

to have a $P\{CS \mid R\}$ of at least $P^* = .95$ for any true point in the parameter space. From Table IB of [6] the appropriate c-value is $c(14, 5, .95) = .2911$. The observed values of T_i are 7.12, 2.02, 3.09, 5.05, 7.28 so that $T_{[1]} = 2.02$ and $T_{[1]}/c = 6.94$. Applying the procedure R, we find that the populations in the selected subset are those three that gave rise to the values 2.02, 3.09, and 5.05. At this point the experimenter can assert with confidence level .95 that one of these three populations has the smallest value of θ . In the case of ranking variances of normal populations $\theta_i = \sigma_i^2$ is the population variance of Π_i and $T_i = vs_i^2$ is proportional to the sample variance from Π_i ($i = 1, 2, \dots, 5$).

5. The Probability of a Correct Selection

In this section exact and asymptotic expressions are derived for the probability of a correct selection for general k , v , and any point in the parameter space and also for the infimum of this probability over all points in the parameter space. It will be convenient here to associate v_i degrees of freedom (the v_i need not be equal) with the sample variance s_i^2 from Π_i ($i = 1, 2, \dots, k$) and to allow the v_i to take on even as well as odd integer values; this generalization is discussed in Section 7. Under this general framework Section 5.2 points out that the $P\{CS \mid R\}$ approaches its infimum as the variances approach equality (with $\sigma_{[1]}^2 = \sigma_{[2]}^2$ only in the limit). It should be noted that the infimum is not the same as the value of the $P\{CS \mid R\}$ when the variances are equal (the latter has not been defined but by any reasonable extension of our definition should be unity); mathematically, we can use the configuration with all variances equal provided we "tag" a particular one of the populations and regard it as being the best. In Section 5.3 the $P\{CS \mid R\}$ is shown to be equivalent to the complement of the c.d.f. of a studentized smallest chi-square statistic and this leads to an approximation of the $P\{CS \mid R\}$ based on asymptotic normality.

5.1 Exact Expressions for the Probability of a Correct Selection

Let $s_{(1)}^2$ denote the (unknown) sample variance that is associated with the i^{th} smallest population variance, $\sigma_{[1]}^2$; let $v_{(1)}$ denote the number of degrees of freedom associated with $s_{(1)}^2$. The procedure R described in Section 3 yields a correct selection if and only if the event

$$(5.1) \quad s_{(1)}^2 \leq \frac{1}{c} \min_{\alpha} s_{(\alpha)}^2 \quad (\alpha = 1, 2, \dots, k)$$

occurs. Since $0 < c \leq 1$, the occurrence of the event (5.1) is equivalent to the occurrence of the event

$$(5.2) \quad s_{(1)}^2 \leq \frac{1}{c} \min_{\alpha \neq 1} s_{(\alpha)}^2 \quad (\alpha = 2, 3, \dots, k).$$

Hence, the probability of a correct selection is given by

$$(5.3) \quad \begin{aligned} P \{ \text{CS} \mid R \} &= P \left\{ s_{(1)}^2 \leq \frac{1}{c} \min_{\alpha \neq 1} s_{(\alpha)}^2 \mid (\alpha = 2, 3, \dots, k) \right\} \\ &= P \left\{ \frac{v_{(\alpha)} s_{(\alpha)}^2}{\sigma_{[\alpha]}^2} \geq \frac{c v_{(\alpha)} \sigma_{[1]}^2}{v_{(1)} \sigma_{[\alpha]}^2} \left(\frac{v_{(1)} s_{(1)}^2}{\sigma_{[1]}^2} \right) \mid (\alpha = 2, 3, \dots, k) \right\} \\ &= \int_0^\infty g_v(x) \prod_{\alpha=2}^k \left[1 - G_v \left(\frac{c v_{(\alpha)} \sigma_{[1]}^2}{v_{(1)} \sigma_{[\alpha]}^2} x \right) \right] dx \end{aligned}$$

where $G_v(x)$ and $g_v(x)$ are the chi-square c.d.f. and p.d.f. with v degrees of freedom, respectively; we can also regard $G(x)$ and $g(x)$ as the gamma c.d.f. and p.d.f., respectively.

Suppose in the next to the last expression of (5.3) we make the transformation

$$(5.4) \quad u_1 = \frac{v_{(1)} s_{(1)}^2}{\sigma_{[1]}^2}, \quad u_\alpha = \frac{v_{(\alpha)} \sigma_{[1]}^2 s_{(\alpha)}^2}{v_{(1)} \sigma_{(\alpha)}^2 s_{(1)}^2} \quad (\alpha = 2, 3, \dots, k)$$

then the limits of u_1 are from 0 to ∞ and upon integrating out u_1 , we obtain the $(k-1)$ -fold integral

$$(5.5) \quad P\{CS \mid R\} = \int_0^\infty \frac{c m_{(2)} \sigma_{[1]}^2}{m_{(1)} \sigma_{[2]}^2} \dots \int_0^\infty \frac{c m_{(k)} \sigma_{[1]}^2}{m_{(1)} \sigma_{[k]}^2} \frac{\Gamma(m_1 + \dots + m_k)}{\prod_{\alpha=1}^k \Gamma(m_\alpha)} \frac{\prod_{\alpha=2}^k \left[u_\alpha^{m_\alpha - 1} du_\alpha \right]}{(1 + u_2 + \dots + u_k)^{m_1 + \dots + m_k}}$$

where $m_{(\alpha)} = v_{(\alpha)}/2$.

5.2 Exact Expressions for the Infimum of the $P\{CS \mid R\}$

It follows from both (5.3) and (5.5) that, for fixed c and fixed $m_{(\alpha)}$ (or $v_{(\alpha)}$), the $P\{CS \mid R\}$ depends only on the ratios of the variances and that it approaches its infimum by setting $\sigma_{[k]}^2 = \sigma_{[k-1]}^2 = \dots = \sigma_{[2]}^2$ letting $\sigma_{[2]}^2 \rightarrow \sigma_{[1]}^2$, with equality only in the limit. Hence, we obtain for the case of a common m [letting Inf_Ω denote the infimum over all points in the parameter space $\sigma_i^2 > 0$ ($i = 1, 2, \dots, k$)],

$$(5.6) \quad \begin{aligned} \text{Inf}_\Omega P\{CS \mid R\} &= \int_c^\infty \dots \int_c^\infty \frac{\Gamma(km)}{[\Gamma(m)]^k} \frac{\prod_{\alpha=2}^k [u_\alpha^{m-1} du_\alpha]}{(1 + u_2 + \dots + u_k)^{km}} \\ &= \int_0^\infty \left[1 - G_v(cx) \right]^{k-1} g_v(x) dx. \end{aligned}$$

In the special case when v is an even integer (so that m is an integer) further simplification takes place in the last expression of (5.6) and we can obtain a result in the form of a finite series.

This result is described in the companion paper [6], where it is used to construct tables for the largest c -value satisfying the probability requirement (3.2).

5.3 Approximation to $\inf P \{CS \mid R\}$ Based on Asymptotic Normality

Let $\chi_j^2(\nu)$ ($j = 0, 1, 2, \dots, p$; $p = k-1$) denote k independent chi-square chance variables with a common number of degrees of freedom ν . It follows from (5.2) and the fact that the $P \{CS \mid R\}$ approaches its infimum as the variances approach equality that we can write the basic probability requirement in the form

$$(5.7) \quad P \left\{ \frac{\min_{j=1, \dots, p} \chi_j^2(\nu)}{\chi_0^2(\nu)} \geq c \right\} = P^* .$$

Hence, the determination of c to satisfy (3.2) for all points in the parameter space is equivalent to the determination of a lower percentage points of the Studentized smallest chi-square statistic with ν degrees of freedom for all k chi-squares.

Using the fact that for large ν the statistic $Y_j = \log (\chi_j^2(\nu)/\nu)$ tends to normality as $\nu \rightarrow \infty$ and also the fact that the set of statistics $\{Y_j - Y_0\}$ tends to a joint p -variate normal distribution, we obtain (the details are similar to those used in Section 5 of [2] for the case $t=1$) the approximation

$$(5.8) \quad \inf_{\Omega} P \{CS \mid R\} \cong \int_{-\infty}^{\infty} \left[1 - F(x-d) \right]^{k-1} f(x) dx$$

$$= \int_{-\infty}^{\infty} \left[F(x+d) \right]^{k-1} f(x) dx$$

where $F(x)$ and $f(x)$ are the standard Normal c.d.f. and p.d.f., respectively, and

$$(5.9) \quad d = \sqrt{(\nu-1)/2} \log (1/c).$$

Another expression for (5.8) is

$$(5.10) \quad \inf_{\Omega} P \left\{ CS \mid R \right\} = (k-1) \int_{-\infty}^{\infty} F(x+d) \left[1-F(x) \right]^{k-2} f(x) dx,$$

and this is the form used in [7] where this quantity is extensively tabulated.

For large values of ν an approximate solution to the largest value of c satisfying the probability requirement (3.2) can be obtained by equating the last expression in (5.8) to P^* and using (5.9) to solve for c . Values in d are tabulated in Table I of [1] for $k = 2(1)10$ and many values of P^* [our d corresponds to $\lambda\sqrt{N}$ in his notation] and also in Table AI of [3] for $k = 2(1)51$ for selected values of P^* [our p and P^* correspond to his n and $1-\alpha$, respectively].

To illustrate numerically the closeness of the normal approximation we take a c -value out of Table IC of [6] corresponding to $P^* = .95$, $k = 3$, $\nu = 50$, namely, $c = .5761$, which is based on an exact calculation, setting the last member of (5.6) equal to $P^* = .95$ and solving for c . This is to be compared with a computation based on (5.8) for which we first must compute d , obtaining

$$(5.11) \quad d = \sqrt{\frac{\nu-1}{2}} \log_e (1/c) = \frac{7\sqrt{2}}{2} \left(\frac{1}{.5761} \right) = 2.729.$$

From the table headed $P(1,3)$ in [7], we obtain by linear interpolation.

$P = .9515$, which is slightly greater than the correct value, namely $P^* = .95$.

6. Expected Size of the Selected Subset

For the procedure R the size S of the selected subset is a chance variable which can take on only integer values from 1 to k , inclusive. For any fixed values of v , k , and P^* , the expected size of the selected subset is a function of the true configuration $\vec{\sigma}^2 = \{\sigma_1^2, \sigma_2^2, \dots, \sigma_k^2\}$ and this function can be regarded as a criterion of the efficiency of any procedure which satisfies the basic probability requirement (3.2). In analogy with power function considerations, one secondary problem is to find the smallest common sample size n necessary to control $E\{S\}$ at some preassigned level for a particular alternative in the parameter space; alternatively, we may wish to control the maximum of $E\{S\}$ over all parameter points in the subset $\Omega(\delta)$ of Ω given by $\sigma_{[i]}^2 \geq \delta \sigma_{[1]}^2$ ($i = 2, 3, \dots, k$) with $\delta > 1$.

6.1 Exact Expression for the Expected Size

Let Y_i denote a chance variable which equals 1 if Π_i is included in the selected subset and equals 0 otherwise. Then $S = \sum_{i=1}^k Y_i$ and hence

$$\begin{aligned} (6.1) \quad E\{S\} &= E\left\{\sum_{i=1}^k Y_i\right\} = \sum_{i=1}^k E\{Y_i\} \\ &= \sum_{i=1}^k P\left\{\Pi_i \text{ is included in the selected subset}\right\} \end{aligned}$$

for any values of v , k , P^* , and $\vec{\sigma}^2$ [these symbols are suppressed in (6.1)].

Using an argument similar to the one in Section 5.1 for obtaining the exact $P\{CS\}$ we obtain from (6.1)

$$(6.2) \quad E\{S\} = \sum_{i=1}^k \int_0^{\infty} g_v(x) \prod_{\substack{j=1 \\ j \neq i}}^k \left[1 - G_v(\theta'_{ij} cx)\right] dx$$

where $\theta'_{ij} = \sigma_i^2 / \sigma_j^2$. [It should be noted that if $\theta_{ij} = \sigma_{[i]}^2 / \sigma_{[j]}^2$ is used in (6.2), then the i^{th} term on the right hand side of (6.2) is the probability that the population associated with $\sigma_{[i]}^2$ is included in the selected subset and hence the sum is again $E \left\{ S \right\}$.]

The particular configuration

$$(6.3) \quad \delta \sigma_{[1]}^2 = \sigma_{[\alpha]}^2, \quad (\alpha = 2, 3, \dots, k)$$

is of some special interest and in this case (6.2) reduces to

$$(6.4) \quad E \left\{ S \right\} = \int_0^\infty \left[1 - G_v \left(\frac{cx}{\delta} \right) \right]^{k-1} g_v(x) dx \\ + (k-1) \int_0^\infty \left[1 - G_v(c\delta x) \right] \left[1 - G_v(cx) \right]^{k-2} g_v(x) dx$$

where the first term on the right side of (6.4) is the $P \left\{ CS \right\}$ for the configuration (6.3). In particular if all the variances are equal then $\delta = 1$ and $E \left\{ S \right\}$ equals $k P^*$.

6.2 Maximum Value of $E \left\{ S \right\}$

It will now be shown that the maximum value of $E \left\{ S \right\}$ takes place when all the population variances are equal. If we set the m largest variances ($1 \leq m < k$) equal to a common value σ^2 (say) and define $\theta_{i.} = \sigma_{[i]}^2 / \sigma^2 = \theta_{.i}^{-1}$ then writing $Q = E \left\{ S \mid \sigma_{[k]}^2 = \sigma_{[k-1]}^2 = \dots = \sigma_{[k-m+1]}^2 = \sigma^2 \right\}$, we obtain from the remark following 6.2,

$$(6.5) \quad Q = \sum_{i=1}^{k-m} \int_0^{\infty} \left[1 - G_V(\theta_i, cx) \right]^m g_V(x) \prod_{\substack{j=1 \\ j \neq i}}^{k-m} \left[1 - G_V(\theta_i, cx) \right] dx \\ + m \int_0^{\infty} \left[1 - G_V(cx) \right]^{m-1} g_V(x) \prod_{j=1}^{k-m} \left[1 - G_V(\theta_j, cx) \right] dx .$$

We now show that the right hand member of (6.5) is a non-increasing function of σ^2 for $\frac{1}{k} \leq P^* \leq 1$ (actually the proof shows that it is a strictly decreasing function of σ^2 for $\frac{1}{k} < P^* < 1$). This proves that it is a maximum when $\sigma^2 = \sigma_{[k-m]}^2$ and since this holds for any integer $m < k$, the desired result will follow.

To show that Q is monotonic we differentiate Q with respect to σ^2 and show that the result is negative for $\frac{1}{k} < P^* < 1$ and $\sigma^2 > \sigma_{[k-m]}^2$. Differentiation gives

$$(6.6) \quad \frac{dQ}{d\sigma^2} = mc \sum_{i=1}^{k-m} \left\{ \frac{\sigma_{[i]}^2}{\sigma^4} \int_0^{\infty} x \left[1 - G_V(\theta_i, cx) \right]^{m-1} g_V(\theta_i, cx) g_V(x) \prod_{\substack{j=1 \\ j \neq i}}^{k-m} \left[1 - G_V(\theta_i, cx) \right] dx \right. \\ \left. - \frac{1}{\sigma_{[i]}^2} \int_0^{\infty} x \left[1 - G_V(cx) \right]^{m-1} g_V(\theta_i, cx) g_V(x) \prod_{\substack{j=1 \\ j \neq i}}^{k-m} \left[1 - G_V(\theta_j, cx) \right] dx \right\} .$$

For $P^* < 1$ it is clear that $c > 0$. If we let $x = x' \theta_i$ in the second integral and drop primes then (6.6) becomes

$$(6.7) \quad \frac{dQ}{d\sigma^2} = mc \sum_{i=1}^{k-m} \frac{\sigma_{[i]}^2}{\sigma^2} \int_0^{\infty} x \left[1 - G_V(\theta_i, cx) \right]^{m-1} \prod_{\substack{j=1 \\ j \neq i}}^{k-m} \left[1 - G_V(\theta_j, cx) \right] .$$

$$\left\{ g_V(\theta_i, cx) g_V(x) - g_V(cx) g_V(\theta_i, x) \right\} dx .$$

It suffices to prove that the expression in braces in (6.7) is negative for each $x > 0$. Treating $g_v(x)$ as a gamma density, we obtain

$$(6.8) \quad \log \left\{ g_v(\theta_i \cdot cx) g_v(x) \div g_v(cx) g_v(\theta_i \cdot x) \right\} = -(1-c) (1-\theta_i) \cdot x.$$

Since $c < 1$ for $P^* > \frac{1}{k}$ it follows from (6.8) that for $\sigma^2 > \sigma_{[k-m]}^2$ the expected size of the selected subset is a decreasing function of σ^2 .

This proves the desired result.

It follows that

$$(6.9) \quad \max_{\Omega} E \left\{ S \right\} = k \int_0^{\infty} [1 - G_v(cx)]^{k-1} g_v(x) dx = k P^*$$

which does not depend on the common value of σ^2 .

The above proof also shows that in the subset $\Omega(\delta)$ of the parameter space the function $E \left\{ S \right\}$ takes on its maximum value when (6.3) holds. Hence, the maximum value is given by (6.4). If it is desired to find the smallest value of the common sample size n (equal to v or $v-1$) such that $E \left\{ S \right\} \leq 1+\epsilon$ for all points in $\Omega(\delta)$ ($\epsilon > 0$ and $\delta > 1$ are preassigned) then the solution in v can be obtained by equating the right hand member of (6.4) to $1+\epsilon$.

6.3 An Approximation to $E \left\{ S \right\}$ Based on Asymptotic Normality

Using the fact that $\log \left(s_{(i)}^2 / \sigma_{[i]}^2 \right)$ ($i = 1, 2, \dots, k$) are independently and asymptotically normally distributed with a common mean and a common variance $2/(v-1)$, we now derive the following approximation for $E \left\{ S \right\}$,

$$(6.10) \quad E \left\{ S \right\} = \sum_{i=1}^k \int_{-\infty}^{\infty} \prod_{\substack{j=1 \\ j \neq i}}^k \left[1 - F \left(x + \frac{\log (c \theta_{ij})}{\sqrt{\frac{2}{v-1}}} \right) \right] f(x) dx.$$

For particular configuration (6.3) this reduces to

$$(6.11) \quad E\{S\} = \int_{-\infty}^{\infty} \left[1-F(x-d-\delta_1)\right]^{k-1} f(x) dx \\ + (k-1) \int_{-\infty}^{\infty} \left[1-F(x-d)\right]^{k-2} \left[1-F(x-d+\delta_1)\right] f(x) dx ,$$

where d is given by (5.9) and $\delta_1 = \sqrt{\frac{v-1}{2}} \log \delta$.

7. Generalizations

In this section we consider two different directions for generalizing the problem treated above. One generalization is to allow unequal degrees of freedom associated with each of the sample variances from the k given populations. The other generalization deals with the problem of selecting a subset containing the t "best" populations, i.e., the populations with the t smallest variances, for $t > 1$. For each of these two generalizations, a general conjecture will be made about the worst configuration and each conjecture will be proved for at least one interesting special case. The appropriate value of c to meet a specified probability requirement has not been computed for either of these generalizations.

7.1 Unequal Degrees of Freedom

In order to find the infimum of the probability of a correct selection, we have to consider both changing the values of the unknown population variances and also different possible pairings of the known degrees of freedom v_1, v_2, \dots, v_k with the ordered population variances. The first part, i.e., changing the values of the population variances, was taken into account in Section 5.2 and it follows from there that we need only consider the case of equal variances (or, more precisely, a sequence of configurations with the smallest population variance

different from the others and with all population variances approaching equality). For the second part, it is clear from (5.3) that we are only concerned with at most k of $k!$ possible pairings since (5.3) depends only on which v is associated with the smallest population variance. We now state a general conjecture about these k possible associations, using the notation of Section 5.2.

Conjecture 1: For any fixed set of k values of v (say, v_1, v_2, \dots, v_k) the probability of a correct selection $P\{CS \mid R\}$ is decreased by an interchange of $v_{(1)}$ and $v_{(\alpha)}$ ($\alpha > 1$) if $v_{(1)} < v_{(\alpha)}$.

[This conjecture can be regarded as the result of two separate conjectures; the first is that the $P\{CS \mid R\}$ is decreased by increasing $v_{(1)}$ holding all the remaining v -values fixed and the second is that for each $\alpha > 1$ the $P\{CS \mid R\}$ is decreased by decreasing $v_{(\alpha)}$ holding all the remaining v -values fixed.]

A proof of this conjecture is given at this point for the following special cases:

Case 1: All k statistics s_i^2 have 2 degrees of freedom, except for one of them which has v degrees of freedom, where $v \geq 2$ is an even integer; let $m = v/2$.

Case 2: There are two populations ($k = 2$) and v_1, v_2 are any integers.

Proof for Case 1:

If we associate m with the best population then using (5.6) the probability of a correct selection becomes

$$(7.1) \quad \int_0^\infty \frac{x^{m-1} e^{-x}}{\Gamma(m)} \left[e^{-\frac{cx}{m}} \right]^{k-1} dx = \left[1 + \frac{c(k-1)}{m} \right]^{-m}$$

and if we associate m with any of the other populations then, again using (5.6), the probability of a correct selection becomes

$$(7.2) \quad \int_0^{\infty} [e^{-cx}]^{k-2} e^{-x} \left[e^{-cmx} \sum_{j=0}^{m-1} \frac{(cmx)^j}{j!} \right] dx = \sum_{j=0}^{m-1} \frac{(cm)^j}{[1+c(k+m-2)]^{j+1}}$$

$$= \frac{1 - \left[\frac{cm}{1+c(k+m-2)} \right]^m}{1 + c(k-2)}$$

Hence, to prove that the conjecture holds in this case it must be shown that for any c ($0 \leq c \leq 1$) and any integer k ($k \geq 2$) the last member of (7.1) is not greater than the last member of (7.2). Since equality holds for $c = 0$, it is sufficient to show that for $c > 0$ and any integer $k \geq 2$

$$(7.3) \quad \left[1+c(k-2) \right] \left[1 + \frac{c(k-1)}{m} \right]^{-m} + \left[1 + \frac{1+c(k-2)}{cm} \right]^{-m} \leq 1.$$

It is well known (and easy to show) that the function $\left(1 + \frac{\Delta}{x}\right)^{-x}$, for any fixed $\Delta > 0$, is a strictly decreasing function of x . It follows that each of the terms on the left side of (7.3) must take its maximum value at $m = 1$. Since equality clearly holds in (7.3) for $m = 1$, the result is proved.

Proof for Case 2:

Assume that $v_1 \leq v_2$. It follows from (5.2) that proving the desired result is equivalent to showing that for any c ($0 < c \leq 1$)

$$(7.4) \quad P \left\{ F_{v_1, v_2} \leq \frac{1}{c} \right\} \geq P \left\{ F_{v_2, v_1} \leq \frac{1}{c} \right\}$$

or equivalently that

$$(7.5) \quad P \left\{ F_{v_1, v_2} > \frac{1}{c} \right\} \leq P \left\{ F_{v_2, v_1} > \frac{1}{c} \right\}$$

where F_{v_1, v_2} is the usual (Snedecor) F chance variable with v_1 and v_2 d.f. Since $1/c > 1$, it is sufficient to show that the densities corresponding to the left and right members of (7.5) satisfy the same inequality for each $F > 1$, i.e., that for $F > 1$ and $v_1 \leq v_2$

$$(7.6) \quad \frac{F^{(v_1-2)/2}}{(v_1+v_2)/2} \leq \frac{F^{(v_2-2)/2}}{(v_1+v_2)/2}$$

By algebraic simplification, this can also be written as

$$(7.7) \quad \frac{F^{v'_1-1}}{v'_1} \leq \frac{F^{v'_2-1}}{v'_2}$$

where $v'_1 = v_1/(v_1+v_2)$ and $v'_2 = v_2/(v_1+v_2)$. Since $v'_1 \leq v'_2$, it is sufficient to show that the function $Q(x) = (F^x - F^{-x})/x$ with $F > 1$ is a strictly increasing function of x for $0 < x < 1$. Letting $Q'(x)$ denote $\frac{dQ(x)}{dx}$ and $Q_1(x) = x^2 Q'(x)$, we obtain by differentiation

$$(7.8) \quad Q_1(x) = x (F^x + F^{-x}) \log F - (F^x - F^{-x})$$

$$(7.9) \quad Q_1'(x) = x (F^x - F^{-x}) (\log F)^2 > 0,$$

using the fact that $x > 0$ and $F > 1$ in the last inequality. Since $Q_1(0) = 0$, it follows from (7.9) that $Q_1(x) > 0$ for $0 < x < 1$ and the same result then holds for $Q'(x)$. It follows that $Q(x)$ is strictly increasing for $0 < x < 1$, as was to be shown.

7.2 Selecting a Subset Containing the t Best Populations

The second generalization deals with the problem of selecting a subset of the k given populations containing the populations with the t smallest variances (with fixed $t \leq k$). In this generalization we will restrict our attention to the case of equal degrees of freedom, for the sake of simplicity. The proposed procedure for this problem is to order the sample variances, (let $s_{[i]}^2$ denote the i^{th} smallest and let s_i^2 denote the one from population Π_i as in Section 3), and put Π_i in the selected subset if and only if

$$(7.10) \quad s_i^2 \leq \frac{s_{[t]}^2}{c}$$

where the constant c with $0 < c < 1$ is evaluated in the "worst configuration" to satisfy a prescribed probability requirement. Then the selected subset must contain at least t populations and if the sample size is too small and/or the population variances are too "close together" then the selected subset may even contain all k populations.

To solve this problem one must first ascertain the worst configuration, i.e., find a sequence of vectors (whose components are population variances) for which the probability of a correct selection approaches its infimum. The above is not easy, principally because of the difficulty in obtaining a simple expression for the probability of a correct selection. We now state a general conjecture about the worst configuration.

Conjecture 2: For any fixed t ($1 \leq t < k$), the $P\left\{CS \mid R_t\right\}$ approaches its infimum if the $k-s+1$ largest population variances ($1 \leq s \leq t$) approach equality (with equality of $\sigma_{[t]}^2$ and $\sigma_{[t+1]}^2$ only in the limit) and the remaining $s-1$ population variances approach (or are set equal to) zero. The integer s takes all values between 1 and t and is a non-decreasing function of the specified probability P^* ; in particular, for P^* close to unity we set $t-1$ population variances equal to zero and the tables for the case $t = 1$ (with the k -value taken to be $k-t+1$) can then be used to determine the appropriate c -value.

It is interesting to note that for $t = 1$ we must have $s = 1$ and the conjecture is consistent with the results of Section 5.2. A proof of the conjecture will be given for the special case $k = 3$, $t = 2$, and a common $v = 2$. In order to prove that the conjecture holds for this special case, it is necessary to obtain an expression for the probability of a correct selection, which we now derive for general t .

We define the integer chance variable T' ($1 \leq T' \leq k$) by saying that $T' = t'$ whenever $s_{(t')}^2 \equiv s_{[t]}^2$, i.e., when $s_{[t]}^2$ comes from the population with variance $\sigma_{[t']}^2$. We define the integer chance variable R ($0 \leq R \leq k-t$) for fixed c by saying that $R = r$ whenever

$$(7.11) \quad s_{[t+r]}^2 < \frac{s_{[t]}^2}{c} < s_{[t+r+1]}^2$$

where $s_{[k+1]}^2 = \infty$ and we disregard all possibilities of equality. Then $s_{[t]}^2$ and $s_{[t]}^2/c$ divide the $k-1$ remaining sample variances into 3 sets S_1, S_2, S_3 of sizes $t-1, r$, and $k-t-r$, respectively. Let $P_1\left\{(S_{1\alpha'}, S_{2\alpha'}, S_{3\alpha'}), t', r; c\right\}$ denote the probability of a correct selection when $T' = t' \leq t$, $R = r$ and $(S_{1\alpha'}, S_{2\alpha'}, S_{3\alpha'})$

represents a partition of the set of $t-1$ "smaller" sample variances and $k-t$ "larger" sample variances that corresponds to a correct decision; let $P\{S_{1\beta}, S_{2\beta}, S_{3\beta} ; t', r; c\}$ denote the same for $T' = t' > t$. Then the probability of a correct selection (letting R_t denote the proposed procedure) can be written as

$$(7.12) \quad P\{CS | R_t\} = \sum_{t'=1}^t \sum_{r=0}^{k-t} \frac{\binom{k-t}{r} \binom{t+r-1}{t-1}}{\alpha=1} P_1\{(S_{1\alpha}, S_{2\alpha}, S_{3\alpha}), t', r; c\} \\ + \sum_{t'=t+1}^k \sum_{r=1}^{k-t} \frac{\binom{k-t-1}{r-1} \binom{t+r-1}{t-1}}{\beta=1} P_2\{(S_{1\beta}, S_{2\beta}, S_{3\beta}), t', r; c\}.$$

If $S_{1\alpha}$ consists of the sample variances $s_{(i)}^2$ for $i = (i_1, i_2, \dots, i_{t-1})$, $S_{2\alpha}$ consists of the $s_{(i)}^2$ for $i = (i_{t+1}, i_{t+2}, \dots, i_{t+r})$ and $S_{3\alpha}$ consists of the $s_{(i)}^2$ for $i = (i_{t+r+1}, i_{t+r+2}, \dots, i_k)$ and if we let $i_t = t'$ then we can write

$$(7.13) \quad P_1\{(S_{1\alpha}, S_{2\alpha}, S_{3\alpha}), t', r; c\} = \int_0^\infty \prod_{j=1}^{t-1} \left[G_v \left(\frac{\sigma_{[t']}^2}{\sigma_{[i_j]}^2} x \right) \right].$$

$$\prod_{j=t+1}^{t+r} \left[G_v \left(\frac{\sigma_{[t']}^2}{\sigma_{[i_j]}^2} \frac{x}{c} \right) - G_v \left(\frac{\sigma_{[t']}^2}{\sigma_{[i_j]}^2} x \right) \right] \prod_{j=t+r+1}^k \left[1 - G_v \left(\frac{\sigma_{[t']}^2}{\sigma_{[i_j]}^2} x \right) \right] g_v(x) dx$$

where $G_v(x)$ and $g_v(x)$ are the chi-square (χ^2) c.d.f. and density, respectively, for v degrees of freedom. The same expression (7.13) also holds for

$P_2 \{ (s_{1\beta}, s_{2\beta}, s_{3\beta}), t', r; c \}$ except that $t' > t$ and hence we now consider partitions of the set of t "small" variances and $k-t-1$ "large" variances that correspond to a correct selection. Some simplification in the first of the two expressions of (7.12) can be obtained by considering separately those cases in which the t best populations yield the t smallest sample variances. The first expression of (7.12) then breaks up into the two parts

$$(7.14) \quad \sum_{t'=1}^t \int_0^{\infty} \prod_{\substack{i=1 \\ i \neq t'}}^t \left[G_v \left(\frac{\sigma_{[t']^2}^2}{\sigma_{[i]}^2} x \right) \right] \prod_{i=t+1}^k \left[1 - G_v \left(\frac{\sigma_{[t']^2}^2}{\sigma_{[i]}^2} x \right) \right] g_v(x) dx$$

$$+ \sum_{t'=1}^t \sum_{r=1}^{k-t} \frac{\binom{k-t}{r} \left[\binom{t+r-1}{t-1} - 1 \right]}{\sum_{\alpha=1}^{\binom{k-t}{r} \left[\binom{t+r-1}{t-1} - 1 \right]}} P_1 \left\{ (s_{1\alpha}, s_{2\alpha}, s_{3\alpha}), t', r; c \right\}$$

where the first part does not depend on c or r and in the second part the index α sums over all the remaining partitions in which at least one ($r > 0$) of $k-t$ populations with larger population variances is associated with one of the $t-1$ smallest observed sample variances.

It should also be noted that if the conjecture is correct then we can set all σ^2 's equal in (7.13) and in (7.14) to obtain a certain amount of further simplification.

Proof of Conjecture for Special Case $k=3$, $t=2$, and Common $v=2$

Suppose we consider the special case $k=3$, $t=2$, and a common v . Then by (7.13) and (7.14) we have

$$\begin{aligned}
 (7.15) \quad P \{CS \mid R_2\} &= P \{s_{(1)}^2 < s_{(2)}^2 < s_{(3)}^2\} + P \{s_{(2)}^2 < s_{(1)}^2 < s_{(3)}^2\} \\
 &+ P \{s_{(3)}^2 < s_{(1)}^2 < s_{(2)}^2 < \frac{1}{c} s_{(1)}^2\} + P \{s_{(3)}^2 < s_{(2)}^2 < s_{(1)}^2 < \frac{1}{c} s_{(2)}^2\} \\
 &+ P \{s_{(1)}^2 < s_{(3)}^2 < s_{(2)}^2 < \frac{1}{c} s_{(3)}^2\} + P \{s_{(2)}^2 < s_{(3)}^2 < s_{(1)}^2 < \frac{1}{c} s_{(3)}^2\}
 \end{aligned}$$

where the 1st, 2nd, and 3rd pair of terms corresponds to the 1st part of (7.14), the 2nd part of (7.14) and the 2nd part of (7.12), respectively. Note that there are symmetries present between the odd and even terms and also between the second and third pair of terms.

Since t is so close to k we gain a little extra simplification by taking complements and writing

$$\begin{aligned}
 (7.16) \quad P \{CS \mid R_2\} &= 1 - \left[P \left\{ s_{(3)}^2 < s_{(1)}^2 < \frac{s_{(1)}^2}{c} < s_{(2)}^2 \right\} + P \left\{ s_{(3)}^2 < s_{(2)}^2 < \frac{s_{(2)}^2}{c} < s_{(1)}^2 \right\} \right. \\
 &\quad \left. - \left[P \left\{ s_{(1)}^2 < s_{(3)}^2 < \frac{s_{(3)}^2}{c} < s_{(2)}^2 \right\} + P \left\{ s_{(2)}^2 < s_{(3)}^2 < \frac{s_{(3)}^2}{c} < s_{(1)}^2 \right\} \right] \right]
 \end{aligned}$$

For $v=2$ this is easily computed and we obtain

$$\begin{aligned}
 (7.17) \quad P \{CS \mid R_2\} &= 1 - \frac{cx}{1+cx} - \frac{c}{x+c} + \frac{c(1+xy)}{c+y+cx} - \frac{c}{c+y} + \frac{c(1+y)}{c+xy+cy} - \frac{c}{c+xy} \\
 &= \frac{x(1-c^2)}{(c+x)(1+cx)} + \frac{cxy^2}{(c+y)(c+y+cx)} + \frac{cxy^2}{(c+xy)(c+xy+cy)}
 \end{aligned}$$

where $x = \sigma_{[2]}^2/\sigma_{[1]}^2$ and $y = \sigma_{[3]}^2/\sigma_{[2]}^2$ are both equal to or greater than unity and $0 < c \leq 1$. Since the first term above is independent of y and each of the last two are strictly increasing in y it follows that for any x we obtain a minimum by setting $y = 1$, obtaining

$$(7.18) \quad P(x \mid y = 1) = \frac{x(1-c^2)}{(c+x)(1+cx)} + \frac{cx}{(c+1)(c+1+cx)} + \frac{cx}{(c+x)(2c+x)},$$

where the symbol $P(x \mid y = 1)$ is used to denote $P\left\{CS \mid R_2\right\}$ for $y = 1$.

We now note that at the extreme values of x we have

$$(7.19) \quad P(1 \mid y = 1) = \frac{3c+1 - 2c^2}{(1+c)(2c+1)}$$

$$(7.20) \quad P(\infty \mid y = 1) = \frac{1}{1+c}$$

and it is easily verified that for $c = \frac{1}{2}$ these two functions are equal with common value $\frac{2}{3}$ and that

$$(7.21) \quad P(\infty \mid y = 1) \begin{matrix} < \\ > \end{matrix} P(1 \mid y = 1) \text{ for } c \begin{matrix} < \\ > \end{matrix} \frac{1}{2}.$$

We now show that for $c \geq \frac{1}{2}$ the minimum for all $x \geq 1$ is given by (7.19) and for $c \leq \frac{1}{2}$ the minimum for all $x \geq 1$ is given by (7.20). In the first part it is sufficient to show that for $x \geq 1$ the difference between (7.18) and (7.19) is non-negative. After some tedious algebra, this difference can be written as

$$(7.22) \quad \frac{c^2(x-1) \left\{ \left[2(c+1)^2 + 4c(1-c^2) \right] + 2x(1+4c^2+2c^3) + x^2(2c-1)(3c^2+2c+2) + x^3c(2c-1) \right\}}{(c+x)(1+cx)(2c+x)(c+1+cx)(c+1)(2c+1)} \geq 0$$

and this proves the first part. For the second part we wish to show that the difference between (7.18) and (7.20) is non-negative and, after simplification, this reduces to showing that for $x \geq 1$ and $0 < c \leq \frac{1}{2}$

$$(7.23) \quad x^2(1-2c^3) + xc(3-4c-2c^2) - 2c \geq 0.$$

To show (7.23), it is sufficient to show for any c with $0 < c \leq \frac{1}{2}$ that the roots of the quadratic in x must be less than or equal to unity. We need only consider the larger of the two roots, which has a plus sign in front of the radical since $1-2c^3 > 0$ for $0 < c \leq \frac{1}{2}$. After simplification, the inequality to be shown takes the form

$$(7.24) \quad (1+c) (1+2c) (1-2c^3) (1-2c) \geq 0$$

which is true for $0 < c \leq \frac{1}{2}$; this proves the second part.

In conclusion, it has been shown that if the specified P^* is between $\frac{1}{3}$ and $\frac{2}{3}$ then we set the right hand member of (7.19) [worst case is $\sigma_{[1]}^2 = \sigma_{[2]}^2 = \sigma_{[3]}^2$] equal to P^* and solve the resulting quadratic, obtaining for the positive root the appropriate value of c ,

$$(7.25) \quad c = \frac{3(1-P^*) + \sqrt{(1-P^*)(17-P^*)}}{4(1+P^*)}$$

If the specified P^* is between $\frac{2}{3}$ and 1 then we use (7.20) [worst case is $\sigma_{[1]}^2 = 0, \sigma_{[2]}^2 = \sigma_{[3]}^2$] in a similar manner and obtain for the appropriate c -value the simple result,

$$(7.26) \quad c = \frac{1-P^*}{P^*}$$

It is of some interest to note that if (7.25) is incorrectly used for P^* between $\frac{2}{3}$ and 1 then the maximum error attained in the guaranteed value of P^* is easily shown to be $7-4\sqrt{3} = .07$ to two decimal places.

REFERENCES

- [1] Bechhofer, R. E., "A Single-Sample Multiple Decision Procedure for Ranking Means of Normal Populations with Known Variances," Ann. Math. Stat., Vol. 25 (1954), pp. 16-29.
- [2] Bechhofer, R. E. and Sobel, M., "A Single-Sample Multiple Decision Procedure for Ranking Variances of Normal Populations," Ann. Math. Stat. Vol 25 (1954), pp. 273-289.
- [3] Gupta S. S., "On a Decision Rule for a Problem in Ranking Means," Mimeograph Series No. 150, May, 1956, Institute of Statistics, University of North Carolina, Chapel Hill, N.C.
- [4] Gupta, S. S. and Sobel, M., "On a Statistic Which Arises in Selection and Ranking Problems," Ann. Math. Stat. Vol. (28), 1957, pp. 957-967.
- [5] Gupta, S. S. and Sobel, M., "Selecting a Subset Containing the Best of Several Binomial Populations," Chapter XX, pp. 224-248 in Contributions to Probability and Statistics, Stanford University Press, 1960.
- [6] Gupta, S. S. and Sobel, M., "The Distribution and Percentage Points of the Smallest of Several Correlated F Statistics," to be submitted for publication.
- [7] National Bureau of Standards, Tables computed for R. E. Bechhofer sponsored by ONR and distributed by the Chemical Corps. Engineering Agency, Army Chemical Center, Maryland, in the pamphlet, "Probabilities Associated with Order Statistics in Samples from Two Normal Populations with Equal Variance," by D. Teichroew (1955).